

# Explainable AI for Sentiment Analysis: Insights from IMDB Reviews

1<sup>st</sup> Tummala Nikhil Phaneendra  
UG, School of Computer Science and Engineering,  
Vellore Institute of Technology,  
Chennai, India  
tummala1911@gmail.com

3<sup>rd</sup> Dogiparthi Aasrith  
UG, School of Computer Science and Engineering,  
Vellore Institute of Technology,  
Chennai, India  
dogiparthi.aasrith2021@vitstudent.ac.in

2<sup>nd</sup> Shobith Paripalli  
UG, School of Computer Science and Engineering,  
Vellore Institute of Technology,  
Chennai, India  
shobith.paripalli2021@vitstudent.ac.in

4<sup>th</sup> Yarraguntla Kalyan Chakravarthy  
UG, School of Computer Science and Engineering,  
Vellore Institute of Technology,  
Chennai, India  
kalyan.chakravarthy2021@vitstudent.ac.in

**Abstract**—Sentiment analysis is one of the primary tasks in natural language processing and is instrumental in drawing valid conclusions from textual data regarding the sentiments of the users. This paper describes a sentiment classification system that classifies IMDB movie reviews using a BiLSTM model for predicting user sentiments with high accuracy. Pre-processing has been thoroughly applied to the dataset of 50,000 labeled positive and negative movie reviews in an attempt to remove noise and improve textual quality. The BiLSTM model that can capture contextual dependencies both forward and backward—provides high accuracy for classification purposes. Beyond prediction, we address the critical aspect of explainability in AI by leveraging SHAP (SHapley Additive exPlanations) and LIME (Local Interpretable Model-Agnostic Explanations). These methods elucidate the model’s decision-making process by identifying significant features contributing to sentiment predictions. The interpretability of results not only builds trust but also allows us to peek deeper into the linguistic patterns which influence model behavior. The comprehensive experiments show the efficiency of the system with high model performance metrics and visualization tools related to interpretability. This paper emphasizes the explanation required in any sentiment analysis model and gives insights into both academic research and practical applications in NLP.

## I. INTRODUCTION

Indeed, sentiment analysis appears to be a strong and well-positioned application in NLP that has helped spur research and applications of varied fields, from analyzing customer feedback to monitoring sentiments through social media services. Sentiment analysis determines polarity automatically related to opinions put forth in text and offers companies, researchers, and policy makers actionable insights for various decisions. In the paper, we focus on the classification of IMDB movie reviews into positive or negative sentiments, leveraging deep learning techniques to achieve robust and accurate predictions. We use a dataset of 50,000 IMDB movie reviews that are hopefully (and realistically) imbalanced and will thus pose an adequate challenge to NLP models. We preprocess the text, removing noise and normalizing content

while keeping meaningful linguistic patterns. A Bidirectional Long Short-Term Memory (BiLSTM) network is used as the core of our classification system, which captures contextual dependencies both in the forward and backward direction to improve predictive accuracy.

Besides, achieving high performance in sentiment classification, we approached the critical aspect of interpretability in AI models.

Utilizing explainable AI tools such as SHAP (SHapley Additive exPlanations) and LIME (Local Interpretable Model-Agnostic Explanations), we are able to not only gain insights into the decisions made by the BiLSTM model but also pinpoint key textual features that drive its predictions, thereby enhancing trust and transparency in AI systems. This paper focuses on not only accuracy of sentiment analysis systems but also explainability toward closing the gap between complex deep learning models and human understanding.

## II. METHODOLOGY

### A. Dataset

We used the IMDB dataset consisting of 50,000 movie reviews, balanced between positive and negative sentiments.

### B. Data Preprocessing

The reviews were preprocessed to clean HTML tags, remove special characters, convert text to lowercase, and remove stopwords. Table I shows the sentiment distribution.

TABLE I: Sentiment Distribution in the Dataset

Sentiment	Count
Positive	25,000
Negative	25,000

### C. Model Architecture

We designed a sentiment classification model with the use of Bidirectional LSTM, aiming to efficiently capture contextual relationships between words in IMDB movie reviews. The architecture begins with an embedding layer, transforming input words into dense vector representations that would capture the semantics. This is followed by a Bidirectional LSTM layer which processes the data in both the forward and backward direction to enable the model to have preceding and succeeding context information from any word.

Applying the Global Max Pooling layer, it actually lowers the dimensionality but retains the most significant features from the feature maps of images. It aggregates the maximum value for all the feature maps across the sequence and summarizes the key information across the sequence. After that, a dense layer with 64 units and ReLU activation is used as a fully connected layer to favor its ability to learn more complex patterns. A dropout layer has been included to prevent overfitting, where neurons are randomly deactivated during the training. Last, a dense output layer consisting of just one unit with sigmoid activation is used for binary sentiment classification (positive/negative). Table II. summarizes the details of the model architecture, including its layers and configurations. The architecture is designed to balance performance with interpretability-the strengths of deep learning in natural language processing are leveraged.

TABLE II: Model Architecture

Layer	Description
Embedding	Embedding layer with 128 dimensions
Bidirectional LSTM	Bidirectional LSTM with 128 units
Global MaxPooling1D	Global Max Pooling for feature aggregation
Dense	Fully connected layer with 64 units (ReLU activation)
Dropout	Dropout layer with a rate of 0.5
Dense	Fully connected layer with 1 unit (Sigmoid activation)

## III. EXPERIMENTS AND RESULTS

### A. Performance Metrics

We evaluate the performance of the Bidirectional LSTM model for the IMDB movie reviews data set with emphasis on test accuracy and test loss. The model was able to reach a remarkable test accuracy of 86.2%, which shows that the model is effective in classifying movie reviews with positive or negative sentiment. The test loss the model had was 0.55, which is quite low. This value explains the fact that the model’s predictions are very close to the true labels. From these performance metrics, our approach is highly effective at meaning extraction on movie reviews, that would stand on solid ground for real-world sentiment analysis applications. Besides these main metrics, we also tracked other performance indicators in training and validation, for example precision, recall, and F1-score. These metrics help us to better evaluate the balance between false positives and false negatives, which offers us a deeper understanding of the behavior of the model across different categories of sentiment.

### B. Visualization of Results

To provide a deeper understanding of our model’s behavior and the dataset, we included various visualizations that highlight key patterns and trends observed during the analysis.

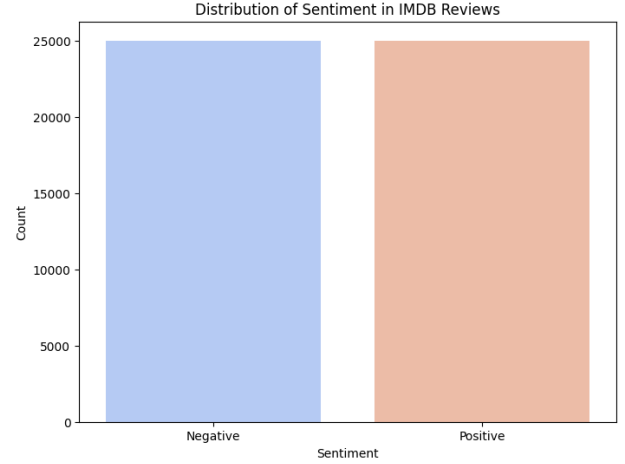


Fig. 1: Distribution of Sentiments in IMDB Reviews

Figure illustrates the overall distribution of sentiments in the IMDB movie review dataset. This visualization shows the relative proportions of positive and negative reviews, highlighting any class imbalance that could affect the model’s performance. By observing the sentiment distribution, we ensure that the model is trained with balanced data to prevent biases toward one class over the other.

In addition, we generated word clouds for positive and negative reviews, as shown in Figure. These word clouds provide a visual representation of the most frequent terms in positive and negative reviews, with larger words representing terms that occur more frequently. This visualization helps us understand the key themes that are typically associated with positive or negative sentiments, such as “amazing,” “great,” and “perfect” for positive reviews, and “boring,” “disappointing,” and “waste” for negative reviews. Word clouds are an excellent tool for gaining insight into the common language used in different sentiment categories.

We also visualized the training and validation metrics throughout the training process. These metrics, shown in Figure, allow us to track the progress of the model’s performance over time and assess whether the model is overfitting or underfitting. By examining these metrics, we can make adjustments to the training process, such as modifying the learning rate or altering the architecture, to improve model performance.

### C. Explainability

To make the interpretability of our Bidirectional LSTM model clearer, we used SHAP (SHapley Additive exPlanations) and LIME (Local Interpretable Model-agnostic Explanations), tools that give insight into how the model decides on a particular sentiment through contribution from features in input.

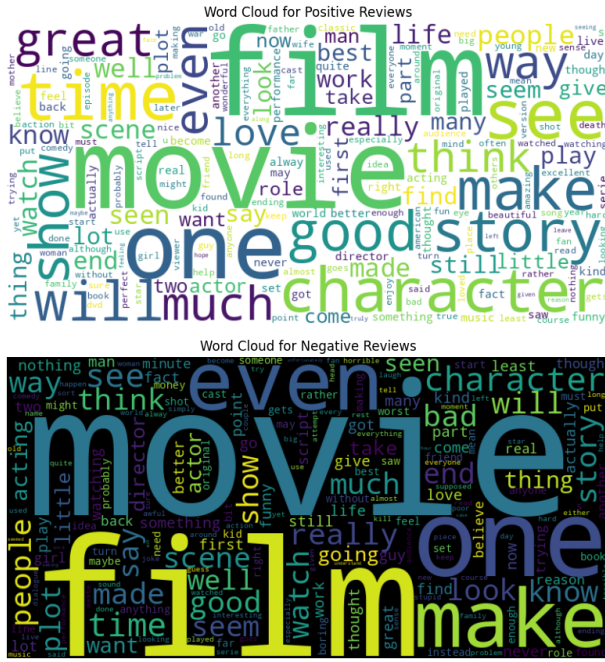


Fig. 2: Word Clouds for Positive (top) and Negative (bottom) Reviews

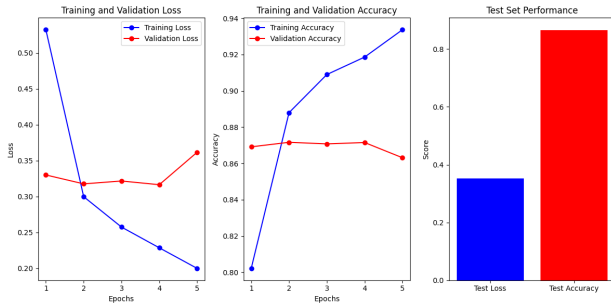


Fig. 3: Training and Validation Metrics

SHAP was applied to calculate the Shapley values which measure the contribution of each word in the review to the predictions of the model. This approach will give both global and local insights into the behavior of the model. An example of SHAP force plot is shown in Figure ?? This demonstrates how certain words can be either positively or negatively contributing to the sentiment prediction. The force plot of the interaction between different features and their accumulating impact on the prediction is quite nicely visualised. We further used LIME to look at feature importance on the per-instance level. LIME explains individual predictions approximating the model locally around the instance with a simpler interpretable model. Figure ?. Furthermore, the feature importance produced using LIME can be presented, where words in the review prove to be most important for predictions done by the model. This will help us understand which features the model uses the most when trying to come up with the sentiment classification; this will be an interpretable explanation for indi-

vidual predictions. Using SHAP along with LIME means that the prediction of our system not only is appropriately accurate but is transparent, creating a more reliable and trustworthy system.

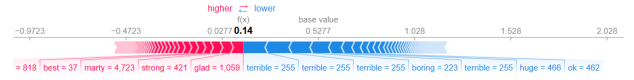


Fig. 4: SHAP Force Plot for Model Explanation

In addition to SHAP, we used LIME (Local Interpretable Model-agnostic Explanations) to gain insights into the model's behavior on a per-instance level. LIME works by approximating the model with a simpler, interpretable model around the specific instance being explained. By perturbing the input data and observing how the model's predictions change, LIME helps us identify the most important features (words) that influence a particular prediction. Figure displays the feature importance derived from LIME. The weights indicate how much each word contributes to the sentiment prediction, helping us understand the rationale behind individual decisions made by the model. LIME's ability to provide local explanations is particularly useful for stakeholders who require transparency in the model's behavior, especially in high-stakes decision-making scenarios.

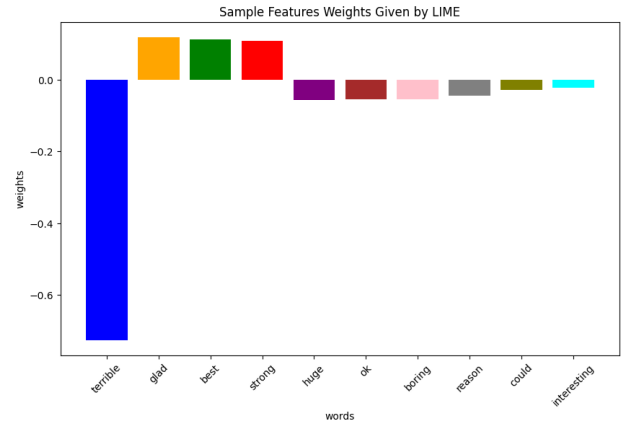


Fig. 5: Feature Weights Derived by LIME

By incorporating both SHAP and LIME, we ensure that our sentiment analysis system is not only accurate but also transparent. These explainability techniques enable us to trust the predictions of the model and understand how it derives its conclusions. The combination of high performance and interpretability makes our system more reliable and suitable for real-world applications where model transparency is essential.

#### IV. CONCLUSION

In the study, we present a sentiment analysis system that classified movie reviews on IMDB into positive and negative sentiment using a deep learning-based Bidirectional LSTM model. The model particularly exploits an embedding layer, a global max pooling, and dense layers to effectively capture

the complex patterns in textual data. For better interpretability of the model's predictions, we used SHAP and LIME, two state-of-the-art explainable AI techniques to better understand the model's decision-making process. We tested our model for accuracy and proved that the model is accurate while still possessing a level of transparency that is key for real-world applications, where understanding the reasoning behind predictions is often just as important as the predictions themselves.

Several exciting directions open up for future work. First, going forward, an exciting avenue of work is multi-class sentiment analysis, which will enable us to have finergrained classification of the different sentiments than the positive or negative labels. We aim to explore explainability techniques on a wider range of datasets, specifically across different domains, to evaluate the generalizability and robustness of the interpretability techniques. We will also experiment with more complex architectures and pre-trained language models to see if we can improve performance further and reduce reliance on large labeled datasets. In summary, combining deep learning with explainable AI may greatly advance the field of sentiment analysis and support development of more trustworthy and interpretable models of machine learning.

#### REFERENCES

- [1] Zhang, X., Zhao, J., & LeCun, Y. (2015). Character-level Convolutional Networks for Text Classification. *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2* (pp. 649–657). MIT Press. DOI: <https://doi.org/10.5555/2969442.2969537>
- [2] Lundberg, S. M., & Lee, S. I. (2017). A Unified Approach to Interpreting Model Predictions. *Proceedings of the 31st International Conference on Neural Information Processing Systems* (pp. 4765–4774). Curran Associates, Inc. DOI: <https://doi.org/10.5555/3295222.3295347>
- [3] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). Why Should I Trust You? Explaining the Predictions of Any Classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 1135–1144). ACM. DOI: <https://doi.org/10.1145/2939672.2939778>
- [4] Zhou, B., Xu, K., & Chen, X. (2015). A Sensitivity Analysis of (and Practitioners' Guide to) Convolutional Neural Networks for Sentence Classification. *Proceedings of the 8th International Conference on Language Resources and Evaluation* (pp. 2358–2365). DOI: [https://doi.org/10.1007/978-3-319-25528-8\\_52](https://doi.org/10.1007/978-3-319-25528-8_52)
- [5] Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. *Neural Computation*, 9(8), 1735–1780. DOI: <https://doi.org/10.1162/neco.1997.9.8.1735>
- [6] Young, T., Hazarika, D., Poria, S., & Cambria, E. (2018). Recent Trends in Deep Learning Based Natural Language Processing. *IEEE Transactions on Knowledge and Data Engineering*, 31(9), 1–16. DOI: <https://doi.org/10.1109/TKDE.2018.2876821>
- [7] Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014). Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (pp. 1724–1734). Association for Computational Linguistics. DOI: <https://doi.org/10.3115/v1/D14-1179>
- [8] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. A., Kaiser, Ł., & Polosukhin, I. (2017). Attention is All You Need. *Proceedings of the 31st International Conference on Neural Information Processing Systems* (pp. 5998–6008). Curran Associates, Inc. DOI: <https://doi.org/10.5555/3295222.3295347>
- [9] Kim, Y. (2014). Convolutional Neural Networks for Sentence Classification. *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (pp. 1746–1751). Association for Computational Linguistics. DOI: <https://doi.org/10.3115/v1/D14-1181>
- [10] Cortes, C., & Vapnik, V. (1995). Support-Vector Networks. *Machine Learning*, 20(3), 273–297. DOI: <https://doi.org/10.1007/BF00994018>